

Here is the theoretical part of Diffusion Model.

1 Langevin SDE

The Langevin SDE has the following form:

$$X_{t+s} = X_t + \nabla \log p_t(x_t)s + \sqrt{2s}\xi \quad (1)$$

where $X_t \in \mathcal{R}^d$, $p_t(x_t) = p(X_t = x_t)$, $\xi \sim N(0, I)$, I is identical matrix of $m \times m$. Our goal is to sample from specific $p(x, t)$.

Theorem 1 *The density of Langevin Diffusion Model converges to $p(x)$ over time. In other words, if $X_t \sim p(x)$, then $X_{t+s} \sim p(x)$ for $\forall s > 0$.*

Proof 1 Let $\mu_t(f) = E[f(X_t)]$. Consider $\mu_{t+\tau}(f) = E[f(X_{t+\tau})]$, as $\tau \rightarrow 0$. Then

$$\begin{aligned} \mu_{t+\tau} &= E \left[f \left(X_t + \nabla \log p_t(x_t) \cdot \tau + \sqrt{2\tau}\xi \right) \right] \\ &= E \left[f(x_t) + \nabla^\top f(x_t) \left(\tau \nabla \log p_t(x_t) + \sqrt{2\tau}\xi \right) \right. \\ &\quad \left. + \frac{1}{2} \left(\nabla^\top \log p_t(x_t) \tau + \sqrt{2\tau}\xi \right) \nabla^2 f(x_t) \nabla \log p_t(x_t) \tau + \sqrt{2\tau}\xi \right] \\ &= E[f(x_t)] + E \left[\tau \nabla^\top f(x_t) \nabla \log p_t(x_t) \right] \\ &\quad + \frac{\tau^2}{2} E \left[\nabla^\top \log p(x_t) \cdot \nabla^2 f(x_t) \cdot \nabla \log p(x_t) \right] + E \left[\tau \xi^\top \nabla^2 f(x_t) \xi \right] \end{aligned} \quad (2)$$

The second term:

$$\begin{aligned} &\tau E \left[\nabla^\top f \nabla \log p_t \right] \\ &= \tau \int \nabla f \cdot \nabla \log p_t p_t dx = \tau \int \nabla f \cdot \nabla p_t dx \\ &= -\tau \int \text{tr}(\nabla^2 f) \cdot p_t dx = -\tau E \left[\text{tr}(\nabla^2 f) \right] \\ &= -\tau E \left[\xi^\top \nabla^2 f \xi \right] \end{aligned} \quad (3)$$

Then

$$\mu_{t+\tau} = E \left[\frac{1}{2} \nabla^\top \log p_t \nabla^2 f \nabla \log p_t \right] \cdot \tau^2 = O(\tau^2) \quad (4)$$

Hence we have $\frac{d}{dt}(\mu_t) = 0$, i.e. $E[\mu_t] = E[\mu_{t+s}]$ for $\forall s > 0$.

Remark 1 We define the density of normal distribution $N(x; \mu, \Sigma)$, and its log-density, gradient of density and score as follows:

$$\begin{cases} N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)} \\ \log N(x; \mu, \Sigma) = -\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu) - \log \left(\sqrt{(2\pi)^d |\Sigma|} \right) \\ \nabla_x N(x; \mu, \Sigma) = N(x; \mu, \Sigma) \Sigma^{-1}(x-\mu) \\ \nabla_x \log N(x; \mu, \Sigma) = -\Sigma^{-1}(x-\mu). \end{cases} \quad (5)$$

Actually, Langevin SDE is not necessary be as above i.e. the diffusion term is not necessary to be $\sqrt{2}$. The reason is to guarantee the stationary distribution of $p_t(x)$. i.e. the term $\frac{\partial p(x,t)}{\partial t} = 0$ in FPK equation. If the diffusion term is $g(t)$, then by FPK equation, we have

$$\nabla_x \cdot (fp - \frac{1}{2}g^2(t)\nabla p) = 0$$

then $f(x, t) = \frac{1}{2}g^2(t) \frac{\nabla_x p(x,t)}{p(x,t)} = \frac{1}{2}g^2(t) \nabla_x \log p(x, t)$.

2 Linear SDE

Then we consider linear SDE having the form:

$$dX_t = (a(t)X_t + b(t))dt + g(t)dW_t \quad (6)$$

By Euler Maruyama method, it can be approximated By

$$\begin{aligned} X_{t+s} &= X_t + (a(t)X_t + b(t))s + g(t)\sqrt{s}\xi \\ &= (1 + a(t)s)X_t + b(t)s + g(t)\sqrt{s}\xi \end{aligned} \quad (7)$$

where $\xi \sim N(0, 1)$. Usually we need to consider the expectation, variance and distribution of x . But the stochastic value of x is dependent of x_0 . Then first we consider

$$\begin{aligned} E[X_{t+s}|X_0] - E[X_t|X_0] &\approx (a(t)E[X_t|X_0] + b(t)s + g(t)\sqrt{s}E[\xi]) \\ &= (a(t)E[X_t|X_0] + b(t))s. \end{aligned} \quad (8)$$

Note $e(t) = E[X_t|X_0]$, then

$$e'(t) = \lim_{s \rightarrow 0} \frac{E[X_{t+s}|X_0] - E[X_t|X_0]}{s} = a(t) \cdot e(t) + b(t). \quad e(0) = X_0. \quad (9)$$

which is an ODE system, having solution

$$e(t) = \left(X_0 + \int_0^t e^{-\int_0^s a(r)dr} b(s)ds \right) \cdot e^{\int_0^t a(s)ds} \quad (10)$$

Therefore

$$\begin{aligned} E[X_t] &= E[E[X_t|X_0]] = E[e(t)] \\ &= \left(E[X_0] + \int_0^t e^{-\int_0^s a(r)dr} b(s)ds \right) e^{\int_0^t a(s)ds} \end{aligned} \quad (11)$$

Similarly, Note $\text{Var}(X_0|X_0) = v(t)$: then $\text{Var}(X_{t+s}|X_0) = (1 + sa(t))^2 \text{Var}(X_t|X_0) + sg^2(t)$. Then

$$\begin{aligned} V'(t) &= \lim_{s \rightarrow 0} \frac{\text{Var}(X_{t+s}|X_0) - \text{Var}(X_t|X_0)}{s} \\ &= [(a^2(t)s + 2a(t))v(t) + g^2(t)]|_{s \rightarrow 0} \\ &= 2a(t)V(t) + g^2(t), \quad V(0) = 0 \end{aligned} \quad (12)$$

Solution is:

$$v(t) = \left(\int_0^t e^{-\int_0^s 2a(r)dr} g^2(s)ds \right) \quad (13)$$

By law of total variance

$$\begin{aligned} \text{Var}(X_t) &= E[X_t^2] - E^2[X_t] = E[E[X_t^2|X_0]] - E^2[X_t] \\ &= E[\text{Var}(X_t|X_0) + E^2[X_t|X_0]] - E^2[X_t] \\ &= E[\text{Var}(X_t|X_0)] + E[E^2[X_t|X_0]] - E^2[E[X_t|X_0]] \\ &= E[\text{Var}(X_t|X_0)] + \text{Var}(E[X_t|X_0]) \end{aligned} \quad (14)$$

then

$$\begin{aligned} \text{Var} &= E[V(t)] + \text{Var}(e(t)) \\ &= \left(\int_0^t e^{-\int_0^s 2a(r)dr} g^2(s)ds \right) e^{\int_0^t 2a(s)ds} + e^{\int_0^t 2a(s)ds} \cdot \text{Var}(X_0). \end{aligned} \quad (15)$$

We have the following theorem which is crucial for diffusion models.

Theorem 2 If $X_{t+s} = (1 + a(t)s)X_t + b(t)s + g(t)\sqrt{s}\xi$
 then $X_t|X_0 \sim N(E[X_t|X_0], \text{Var}(X_t|X_0))$, where $(E[X_t|X_0] = e(t), \text{Var}(X_t|X_0) = V(t))$.

It should be noted that $e(t)$ is related to X_0 and t , while $V(t)$ only depends on t !

Next, we will see how the above formula can be applied to diffusion models. There are three frameworks to build SDEs for diffusion models, VP, VE and sub-VP.

Definition 1 Noise function $\beta(t)$. s.t. $\beta(0) = 0; \beta'(t) \geq 0; \beta(t) \rightarrow \infty$ as $t \rightarrow \infty$

2.1 Variance Preserving (VP) SDE

So if we have diffusion model like:

$$\begin{aligned} X_{t_{i+1}} &= \sqrt{1 - (\beta(t_{i+1}) - \beta(t_i))}X_{t_i} + \sqrt{(\beta(t_{i+1}) - \beta(t_i))}\xi \\ &= \sqrt{1 - \Delta\beta(t_i)}X_{t_i} + \sqrt{\Delta\beta(t_i)}\xi \end{aligned} \quad (16)$$

Then the conditional distribution is given by:

$$q(X_{t_{i+1}}|X_{t_i}) = N(x_{t_{i+1}}; \sqrt{1 - \Delta\beta(t_i)}X_{t_i}, \Delta\beta(t_i)) \quad (17)$$

Then we need to estimate θ drift term f and diffusion term g :

$$\begin{aligned} f(x, t) &= \lim_{h \rightarrow 0} \frac{E[X_{t+h} - X_t | X_t = x]}{h} \\ &= \lim_{h \rightarrow 0} \frac{x\sqrt{1 - \Delta\beta(t)} - x}{h} = -\frac{x}{2}\beta'(t). \\ g(t) &= \sqrt{\lim_{h \rightarrow 0} \frac{N[X_{t+h}|X_t = x]}{h}} = \sqrt{\lim_{h \rightarrow 0} \frac{\beta(t+h) - \beta(t)}{h}} = \sqrt{\beta'(t)} \end{aligned} \quad (18)$$

Then the model can be written as $dx = -\frac{x}{2}\beta'(t)dt + \sqrt{\beta'(t)}dW_t$

Then by Theorem 2 we have

$$\begin{cases} E[X_t|X_0] = X_0 e^{\int_0^t -\frac{1}{2}\beta'(s)ds} = X_0 e^{-\frac{1}{2}\beta(t)} \\ E[X_t] = E[X_0] e^{-\frac{1}{2}\beta(t)} \\ V(X_t|X_0) = \int_0^t e^{\int_0^s \beta'(r)dr} \beta'(s)ds \cdot e^{-\beta(t)} = 1 - e^{-\beta(t)} \\ V(X_t) = 1 - e^{-\beta(t)} + V(X_0) e^{-\beta(t)} = 1 + (V(X_0) - 1) e^{-\beta(t)}. \end{cases} \quad (19)$$

So as $t \rightarrow \infty, \beta(t) \rightarrow \infty$, then $E \rightarrow 0, V \rightarrow 1$, i.e. $X_t|X_0 \sim N(E[X_t|X_0], \text{Var}[X_t|X_0]) \rightarrow N(0, 1)$ as $t \rightarrow \infty$.

2.2 Variance-Exploding SDE

Here is the model: $X_{t+h} = X_t + \sqrt{\Delta\beta(t)}\xi$

Similarly we can compute the $f(x, t) \equiv 0$ and $g(t) = \sqrt{\beta'(t)}$. Hence

$$\begin{cases} E[X_0|X_0] = X_0 \\ E[X_t] = E[X_0] \\ V(X_t|X_0) = \int_0^t e^{\int_0^s 0dr} \beta'(s)ds = \beta(t) \\ V(X_t) = V[X_0] + \beta(t) \end{cases} \quad (20)$$

So the expectation value is constant and the variance is increasing monotonical.

If we rescale X_t as $Y_t = \frac{X_t}{\sqrt{\beta(t)}}$, then $Y_t \rightarrow N(0, 1), t \rightarrow \infty$.

2.3 Sub-VP SPE

Here, we set the drift and diffusion term as

$$\begin{aligned} f(x, t) &= -\frac{1}{2}\beta'(t) \\ g(t) &= \sqrt{\beta'(t)(1 - e^{-2\beta(t)})} \end{aligned} \quad (21)$$

As the same, we can compute that.

$$\begin{cases} E[X_t|X_0] = X_0 e^{-\frac{1}{2}\beta(t)} \\ E[X_t] = E[X_0] e^{-\frac{1}{2}\beta(t)} \\ V(X_t|X_0) = (1 - e^{-\beta(t)})^2 \\ V(X_t) = (1 - e^{-\beta(t)})^2 + V(X_0) e^{-\beta(t)}. \end{cases} \quad (22)$$

We can find out that the variance is always smaller than VP SDE.

Remark 2 To sum up, finally we hope that X_t converges to a normal distribution by choosing different drift and diffusion functions. For generative model, the goal is to sample from a Data distribution p_{data} . We have known that if we set the initial distribution $p_0(x_0) = p(X_0 = x_0) \sim p_{data}$, then after $t = T$, the distribution of X_t is tend to be $N(0, 1)$ under certain conditions.

So the idea is backward: if we sample from $X_T \sim N(0, 1)$, and then run SDE backwards, could we get the initial distribution?

3 Reverse SDE

Assume we have forward SDE: from $X_0 \sim p_0, X_T \sim p_T$,

$$dX_t = f(X_t, t)dt + g(t)dW_t \quad (23)$$

Then we define the reverse SDE as: from $X_T \sim p_T$,

$$d\bar{X}_t = \bar{f}(\bar{X}_t, t)dt + \bar{g}(t)d\bar{W}_t \quad (24)$$

where \bar{W}_t is Brownian Motion runs backward in time, i.e. $\bar{W}_{t-s} - \bar{W}_t$ is independent of \bar{W}_t . We can approximate by EM:

$$\bar{X}_{t-s} - \bar{X}_t = -s\bar{f}(\bar{X}_t, t) + \sqrt{s}\bar{g}(t)\xi \quad (25)$$

So the problem is: If given f, g , are there \bar{f}, \bar{g} s.t. the reverse time diffusion process \bar{X}_t has the same distribution as the forward process X_t ? Yes!

Theorem 3 The reverse SDE with \bar{f}, \bar{g} having the following form has the same distribution as the forward SDE 23:

$$\begin{cases} \bar{f}(x, t) = f(x, t) - g^2(x, t) \frac{\partial}{\partial x} \log p_t(x) \\ \bar{g} = g(t) \end{cases} \quad (26)$$

i.e.

$$d\bar{X}_t = [f(\bar{X}_t, t) - g^2(t) \log p_t(x_t)] dt + g(t)d\bar{W}_t \quad (27)$$

Proof 2 The proof is skipped.

This theorem allows us to learn how to generate samples from p_{data} .

Algorithm 1 :

Step1. Select $f(x, t)$ and $g(t)$ with affine drift coefficients s.t. $X_T \sim N(0, 1)$

Step2. Train a network $s_\theta(x, t) = \frac{\partial}{\partial x} \log p_t(x)$ where $p_t(x) = p(X_t = x)$ is the forward distribution.

Step3. Sample X_T from $N(0, 1)$, then run reverse SDE from T to 0 :

$$\bar{X}_{t-s} = \bar{X}_t + s [g^2(t)s_\theta(\bar{X}_t, t) - f(\bar{X}_t, t)] + \sqrt{s}g(t)\xi \quad (28)$$

4 Loss function

Normally we can define the loss function as follows:

$$\begin{aligned} L_\theta &= \frac{1}{T} \int_0^T \lambda(t) E_{x_0 \sim p_{data}} \left[E_{x_t \sim p_{t|0}(x_t|x_0)} [\|s_\theta(x_t, t) - \nabla_{x_t} \log p_t(x_t)\|^2] \right] dt \\ &= E_{t \sim U(0, T)} \left[\lambda(t) E_{x_0 \sim p_{data}} \left[E_{x_t \sim p_{t|0}(x_t|x_0)} [\|s_\theta(x_t, t) - \nabla_{x_t} \log p_t(x_t)\|^2] \right] \right] \end{aligned} \quad (29)$$

It should be clarified that $p_{t|0}(x_t|x_0) = p(X_t = x_t | X_0 = x_0)$. So

$$p_t(x_t) = \int p_{t|0}(x_t|x_0)p_0(x_0)dx_0 = E_{x_0 \sim p_{data}} [p_{t|0}(x_t|x_0)]$$

. Then $p_t(x) = p(X_t = x)$. $p_{t|0}(x|y) = p(X_t = x | x_0 = y)$.

$$\begin{aligned} \nabla \log p_t(x) &= \frac{1}{p_t(x)} \nabla p_t(x). \\ &= \frac{1}{p_t(x)} \nabla \int p_{t|0}(x|y)p_0(y)dy \\ &= \frac{1}{p_t(x)} \int \nabla p_{t|0}(x|y)p_0(y)dy \\ &= \frac{1}{p_t(x)} \int \frac{\nabla p_{t|0}(x|y)}{p_{t|0}(x|y)} p_0(y) \cdot p_{t|0}(x|y)dy \\ &= \int \nabla_x \log(p_{t|0}(x|y)) \cdot p_0(y)dy \\ &= E_{y \sim p_0(x|y)} [\nabla_x \log(p_{t|0}(x|y))] \end{aligned} \quad (30)$$

Where we have used the following lemma:

Lemma 1

$$E_{x_0 \sim p_0} \left[E_{x_t \sim p_{t|0}(\cdot|x_0)} \left[E_{x'_0 \sim p_{0|t}(\cdot|x_t)} [f(x_t, x'_0)] \right] \right] = E_{x_0 \sim p_0} \left[E_{x_t \sim p_{t|0}(\cdot|x_0)} [f(x_t, x_0)] \right] \quad (31)$$

Proof 3 Easy to prove.

Then we can rewrite the loss function as:

$$\begin{aligned} L_\theta &= E_{t \sim U(0, T)} \left[\lambda(t) E_{x_0 \sim p_{data}} \left[E_{x_t \sim p_{t|0}(x_t|x_0)} \left[\left\| S_\theta(x_t, t) - \frac{\partial}{\partial x_t} \log p_t(x_t) \right\|^2 \right] \right] \right] \\ &\leq E_{t \sim U(0, T)} \left[\lambda(t) E_{x_0 \sim p_{data}} \left[E_{x_t \sim p_{t|0}(x_t|x_0)} \left[E_{y \sim p_{data}} [\|S_\theta(x_t, t) - \nabla_{x_t} \log(p_{t|0}(x_t|y))\|^2] \right] \right] \right] \\ &= E_{t \sim U(0, T)} \left[\lambda(t) E_{x_0 \sim p_{data}} \left[E_{x_t \sim p_{t|0}(x_t|x_0)} [\|S_\theta(x_t, t) - \nabla_{x_t} \log(p_{t|0}(x_t|x_0))\|^2] \right] \right] \end{aligned} \quad (32)$$

Since $p_{t|0}(x_t|x_0) = p(X_t = x_t | X_0 = x_0)$ has been discussed:

$$p_{t|0}(x_t|x_0) \sim N(x_t; E[X_t = x_t | X_0 = x_0], \text{Var}(X_t = x_t | X_0 = x_0)).$$

Then by theorem 2, x can be written as $x = e(t, X_0) + \sqrt{V(t)}\xi$, where $\xi \sim N(0, 1)$, then the score function is:

$$\frac{\partial}{\partial x} \log p_{t|0}(x|x_0) = -\frac{x - E_{t|0}[x|x_0]}{\text{Var}_{t|0}(x|x_0)} = -\frac{x - e(t, X_0)}{V(t)} \sim -N\left(0, \frac{1}{V(t)}\right) \quad (33)$$

So

$$\begin{aligned} L_\theta &= E_{t \sim U(0, T)} \left[\lambda(t) E_{x_0 \sim p_{data}} \left[E_{\xi \sim N(0, 1)} \left[\left\| s_\theta \left(\sqrt{V(t)}\xi + e(t, X_0), t \right) + \frac{\xi}{\sqrt{V(t)}} \right\|^2 \right] \right] \right] \\ &= E_{t \sim U(0, T)} \left[\lambda(t) E_{x_0 \sim p_{data}} \left[\frac{1}{V(t)} E_{\xi \sim N(0, 1)} \left[\left\| \xi_\theta \left(\sqrt{V(t)}\xi + e(t, X_0), t \right) - \xi \right\|^2 \right] \right] \right] \end{aligned} \quad (34)$$

where $\xi_\theta = -\sqrt{V(t)}s_\theta$ is called denoising network.

5 With Classifier

Though we can produce pictures by sampling from normal distribution, we still cannot control what we will generate. What we want to do is something like: "Give me the pictures of number 6", then the model can sample from the normal distribution and do the denoising to generate pics of 6.

Usually, we can do something like: train a model for every class label. This do make the model smaller, but increases number of models. Think about it, when the label is TEXT, it is impossiable to train a model for each sentences.

So, the initial distribution is $p_0(x|y)$ given the label y . Similarly, we will convert the data distribution $p_{data}(x|y)$ to final distribution, normal distribution expected. Then we SDE becomes: $X_t \sim p_t(x|y)$

$$\begin{aligned} P_t(x | y) &= P(X_t = x | y) = \frac{P(y | X_t = x) P(X_t = x)}{P(y)} \\ \Rightarrow \log(P_t(x | y)) &= \log(P(y | X_t = x)) + \log(P(X_t = x)) - \log(P(y)) \\ \Rightarrow \nabla_x \log(P_t(x | y)) &= \nabla_x \log(P(y | X_t = x)) + \nabla_x \log(P(X_t = x)) \end{aligned} \quad (35)$$

We have finished training $\nabla_x \log(p(X_t = x))$ in sampling. Then we need to estimate $\nabla_x \log(p(y | X_t = x))$. This is the conditional protability, we end up with a sharp factor s: $P'(y | X_t = x)$, then:

$$\nabla_x \log(p_t(x | y)) = S \nabla_x \log(p(y | x_t = x)) + \nabla_x \log(p(x_t = x)) \quad (36)$$

Note $\omega_\theta(y | x, t)$ to learn $s \nabla_x \log(p(y | X_t = x))$